

Bayesian Models in Machine Learning

Approximate inference in Bayesian models
Latent Dirichlet Allocation Model

Lukáš Burget



Latent Dirichlet Allocation Model

Set of topic (color) specific word distributions

$$\Phi^T = \begin{bmatrix} \boldsymbol{\varphi}_1^T \\ \boldsymbol{\varphi}_2^T \\ \boldsymbol{\varphi}_3^T \end{bmatrix} = \begin{array}{c} \text{tenis} \\ \text{surfing} \\ \text{software} \\ \text{apple} \\ \text{burger} \\ \dots \end{array} \begin{bmatrix} 0.4 & 0.2 & 0.0 & 0.0 & 0.0 & \dots \\ 0.0 & 0.4 & 0.2 & 0.2 & 0.0 & \dots \\ 0.0 & 0.0 & 0.0 & 0.6 & 0.2 & \dots \end{bmatrix} \begin{array}{c} \text{sports} \\ \text{computers} \\ \text{food} \end{array} \begin{bmatrix} 0.5 \\ 0.0 \\ 0.5 \end{bmatrix} = \boldsymbol{\theta}_d$$

Document specific topic mixture weights

$$\Phi \boldsymbol{\theta}_d = [0.2 \quad 0.1 \quad 0.0 \quad 0.3 \quad 0.1 \quad \dots]$$

$$\mathbf{w}_d = [\text{apple} \quad \text{burger} \quad \text{is} \quad \text{surfing} \quad \text{apple} \quad \text{tenis}]$$

Document specific word distribution.

Mixture of Categorical distributions is again Categorical distribution.

Document, where each word is independently drawn from the document specific distribution

$$w_{jt} \sim \text{Cat}(\Phi \boldsymbol{\theta}_d)$$

The Task

Given a set of documents and chosen number of topics

learn the set of topic specific word distribution

without knowing any topic labels (i.e. topics are latent)

and, for each document, estimate the topic mixture weights.

$$\Phi^T = \begin{bmatrix} \boldsymbol{\varphi}_1^T \\ \boldsymbol{\varphi}_2^T \\ \boldsymbol{\varphi}_3^T \end{bmatrix} = \begin{array}{c} \textit{tenis} \\ \textit{surfing} \\ \textit{software} \\ \textit{apple} \\ \textit{burger} \end{array} \begin{bmatrix} 0.4 & 0.2 & 0.0 & 0.0 & 0.0 & \dots \\ 0.0 & 0.4 & 0.2 & 0.2 & 0.0 & \dots \\ 0.0 & 0.0 & 0.0 & 0.6 & 0.2 & \dots \end{bmatrix} \begin{array}{c} \textit{sports} \\ \textit{computers} \\ \textit{food} \end{array} \begin{bmatrix} 0.5 \\ 0.0 \\ 0.5 \end{bmatrix} = \boldsymbol{\theta}_d$$

$$\Phi \boldsymbol{\theta}_d = [0.2 \quad 0.1 \quad 0.0 \quad 0.3 \quad 0.1 \quad \dots]$$

- $\mathbf{w}_1 = [\textit{apple} \quad \textit{burger} \quad \textit{is} \quad \textit{surfing} \quad \textit{apple} \quad \textit{tenis}]$
- $\mathbf{w}_2 = [\textit{tenis} \quad \textit{and} \quad \textit{surfing} \quad \textit{but} \quad \textit{mainly} \quad \textit{tenis}]$
- $\mathbf{w}_3 = [\textit{surfing} \quad \textit{with} \quad \textit{apple} \quad \textit{software} \quad \textit{about} \quad \textit{tenis}]$
- \vdots
- $\mathbf{w}_D = [\textit{the} \quad \textit{best} \quad \textit{burget} \quad \textit{is} \quad \textit{appe} \quad \textit{burger}]$

This vector can serve as low-dimensional representation of the document (e.g. for topic clustering).
Group of documents for which the same weight dominates are probably on one and the same topic.

LDA assumed generative process

- For each document d , each word w_{dn} is independently drawn from the document specific distribution $\Phi\theta_d$:

for $d = 1..D$

for $n = 1..N_d$

$$w_{dn} \sim \text{Cat}(\Phi\theta_d)$$

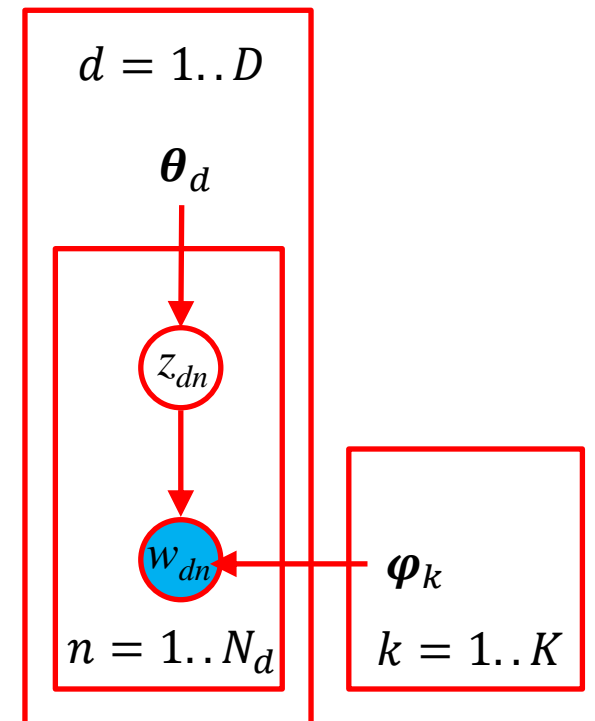
- Or, like for GMM, we choose the mixture component z_{dn} (representing a topic) and sample observation w_{dn} from its distribution.

for $d = 1..D$

for $n = 1..N_d$

$$z_{dn} \sim \text{Cat}(\theta_d)$$

$$w_{dn} \sim \text{Cat}(\varphi_{z_{dn}})$$

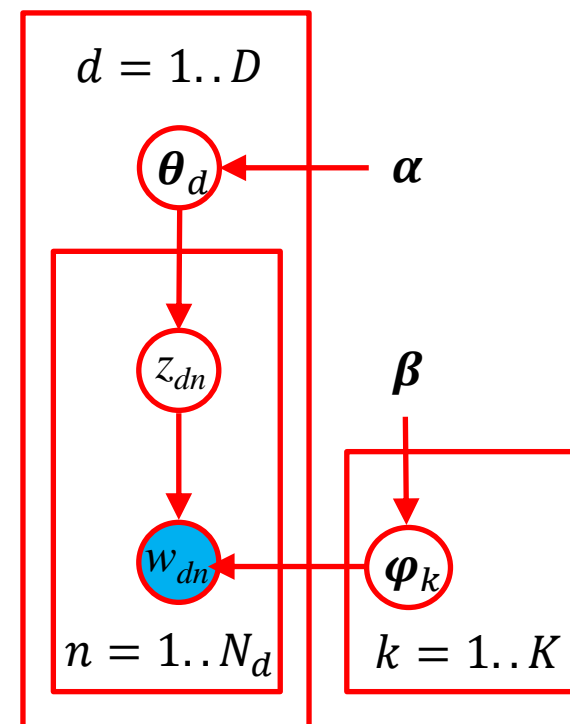


Full Bayesian LDA

- Let's treat parameters of each topic specific distribution φ_k and document specific weights θ_d as random variables with Dirichlet priors: $\varphi_k \sim \text{Dir}(\beta)$ and $\theta_d \sim \text{Dir}(\alpha)$

- The generative process is now:

for $k = 1..K$
 $\varphi_k \sim \text{Dir}(\beta)$
for $d = 1..D$
 $\theta_d \sim \text{Dir}(\alpha)$
for $n = 1..N_d$
 $z_{dn} \sim \text{Cat}(\theta_d)$
 $w_{dn} \sim \text{Cat}(\varphi_{z_{dn}})$



Bayesian LDA model summary

Joint Probability:

$$p(\mathbf{W}, \mathbf{Z}, \Theta, \Phi) = \left(\prod_{k=1}^K p(\boldsymbol{\varphi}_k) \right) \left(\prod_{d=1}^D p(\boldsymbol{\theta}_d) \right) \left(\prod_{d=1}^D \prod_{n=1}^{N_d} P(z_{dn} | \boldsymbol{\theta}_d) \right) \left(\prod_{d=1}^D \prod_{n=1}^{N_d} P(w_{dn} | \Phi, z_{dn}) \right)$$

Variables:

$\Phi = [\boldsymbol{\varphi}_1, \boldsymbol{\varphi}_2, \dots, \boldsymbol{\varphi}_K]$ - topic specific word distributions

$\Theta = [\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_D]$ - document specific topic distributions

$z_{dn} = k$ - denotes that n^{th} word in document d comes from topic k

$w_{dn} = v$ - denotes that n^{th} word in document d is v

Indices:

$d = 1..D$ - (training) document

$k = 1..K$ - topic

$v = 1..V$ - unique word in the vocabulary of size V

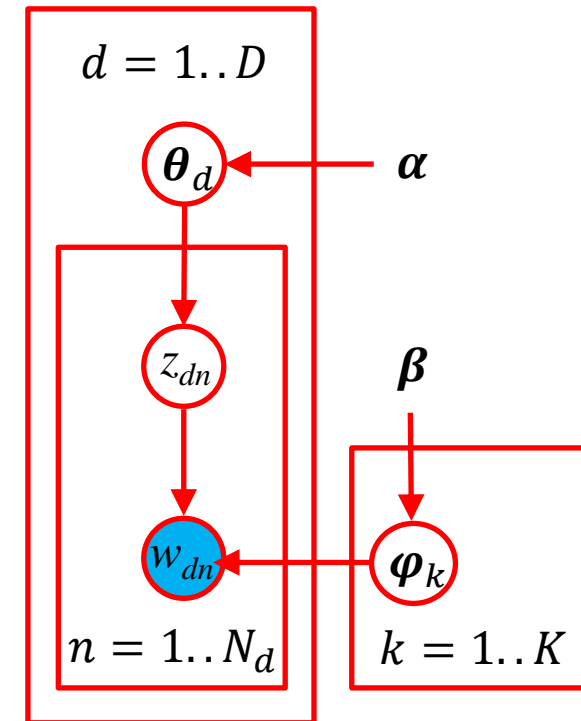
$n = 1..N_d$ - position of word in document d

$p(\boldsymbol{\varphi}_k) = \text{Dir}(\boldsymbol{\varphi}_k | \boldsymbol{\beta}) \propto \prod_{v=1}^V \varphi_{kv}^{\beta_v - 1}$ - prior on parameters $\boldsymbol{\varphi}_k$

$p(\boldsymbol{\theta}_d) = \text{Dir}(\boldsymbol{\theta}_d | \boldsymbol{\alpha}) \propto \prod_{k=1}^K \theta_{dk}^{\alpha_k - 1}$ - prior on document specific topic mixture weights

$P(z_{dn} = k | \boldsymbol{\theta}_d) = \theta_{dk}$ - probability that a word in document d comes from topic k

$P(w_{dn} = v | \Phi, z_{dn} = k) = \varphi_{kv}$ - probability of word v if we know that it comes from topic k



Using counts $C_{d,v}^k$

Joint Probability:

$$p(\mathbf{W}, \mathbf{Z}, \Theta, \Phi) = \left(\prod_{k=1}^K p(\boldsymbol{\varphi}_k) \right) \left(\prod_{d=1}^D p(\boldsymbol{\theta}_d) \right) \left(\prod_{d=1}^D \prod_{n=1}^{N_d} P(z_{dn} | \boldsymbol{\theta}_d) \right) \left(\prod_{d=1}^D \prod_{n=1}^{N_d} P(w_{dn} | \Phi, z_{dn}) \right)$$

Let $C_{d,v}^k$ be the count of words v generated from topic k in document d as assigned by latent variables z_{dn} .

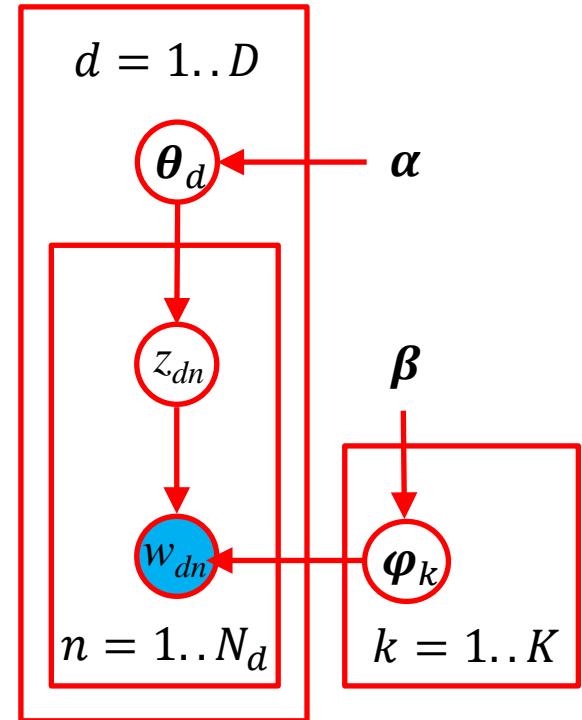
$$C_{d,v}^k = \sum_{n=1}^{N_d} \delta(z_{dn} = k) \delta(w_{dn} = v)$$

$$\prod_{n=1}^{N_d} P(z_{dn} | \boldsymbol{\theta}_d) = \prod_{n=1}^{N_d} \theta_{dz_{dn}} = \prod_{k=1}^K \theta_{dk}^{C_{d,(\cdot)}^k},$$

where $C_{d,(\cdot)}^k = \sum_{v=1}^V C_{d,v}^k$ is the count of all words from topic k in document d

$$\prod_{d=1}^D \prod_{n=1}^{N_d} P(w_{dn} | \Phi, z_{dn}) = \prod_{d=1}^D \prod_{n=1}^{N_d} \varphi_{z_{dn} w_{dn}} = \prod_{k=1}^K \prod_{v=1}^V \varphi_{kv}^{C_{(\cdot),v}^k},$$

where $C_{(\cdot),v}^k = \sum_{d=1}^D C_{d,v}^k$ is the count of words v from topic k in all documents.

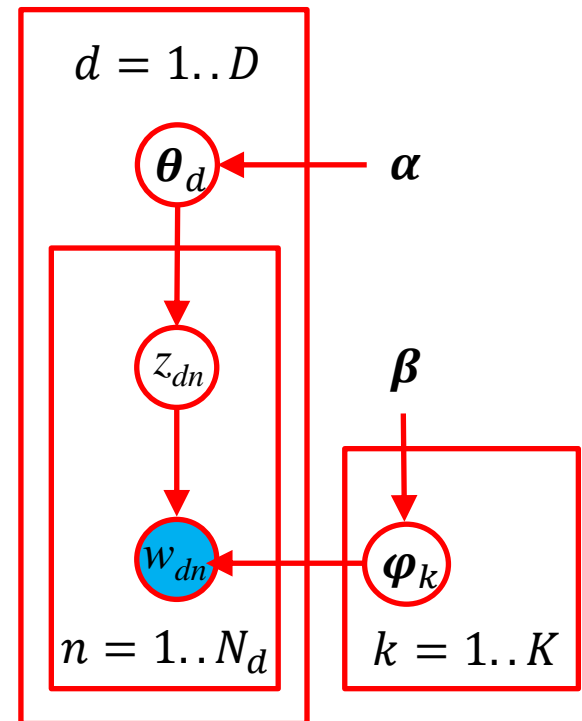


Joint probability using the counts

$$\begin{aligned}
 p(\mathbf{W}, \mathbf{Z}, \boldsymbol{\Theta}, \boldsymbol{\Phi}) &= \left(\prod_{k=1}^K p(\boldsymbol{\varphi}_k) \right) \left(\prod_{d=1}^D p(\boldsymbol{\theta}_d) \right) \left(\prod_{d=1}^D \prod_{n=1}^{N_d} P(z_{dn} | \boldsymbol{\theta}_d) \right) \left(\prod_{d=1}^D \prod_{n=1}^{N_d} P(w_{dn} | \boldsymbol{\Phi}, z_{dn}) \right) \\
 &= \left(\prod_{k=1}^K \frac{\Gamma(\sum_{v=1}^V \beta_v)}{\prod_{v=1}^V \Gamma(\beta_v)} \prod_{v=1}^V \varphi_{kv}^{\beta_v-1} \right) \left(\prod_{d=1}^D \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \theta_{dk}^{\alpha_k-1} \right) \left(\prod_{d=1}^D \prod_{k=1}^K \theta_{dk}^{C_{d,(\cdot)}^k} \right) \left(\prod_{k=1}^K \prod_{v=1}^V \varphi_{kv}^{C_{(\cdot),v}^k} \right) \\
 &\propto \left(\prod_{k=1}^K \prod_{v=1}^V \varphi_{kv}^{\beta_v-1} \right) \left(\prod_{d=1}^D \prod_{k=1}^K \theta_{dk}^{\alpha_k-1} \right) \left(\prod_{d=1}^D \prod_{k=1}^K \theta_{dk}^{C_{d,(\cdot)}^k} \right) \left(\prod_{k=1}^K \prod_{v=1}^V \varphi_{kv}^{C_{(\cdot),v}^k} \right)
 \end{aligned}$$

$$C_{d,v}^k = \sum_{n=1}^{N_d} \delta(z_{dn} = k) \delta(w_{dn} = v)$$

$$C_{d,(\cdot)}^k = \sum_{v=1}^V C_{d,v}^k \quad C_{(\cdot),v}^k = \sum_{d=1}^D C_{d,v}^k$$

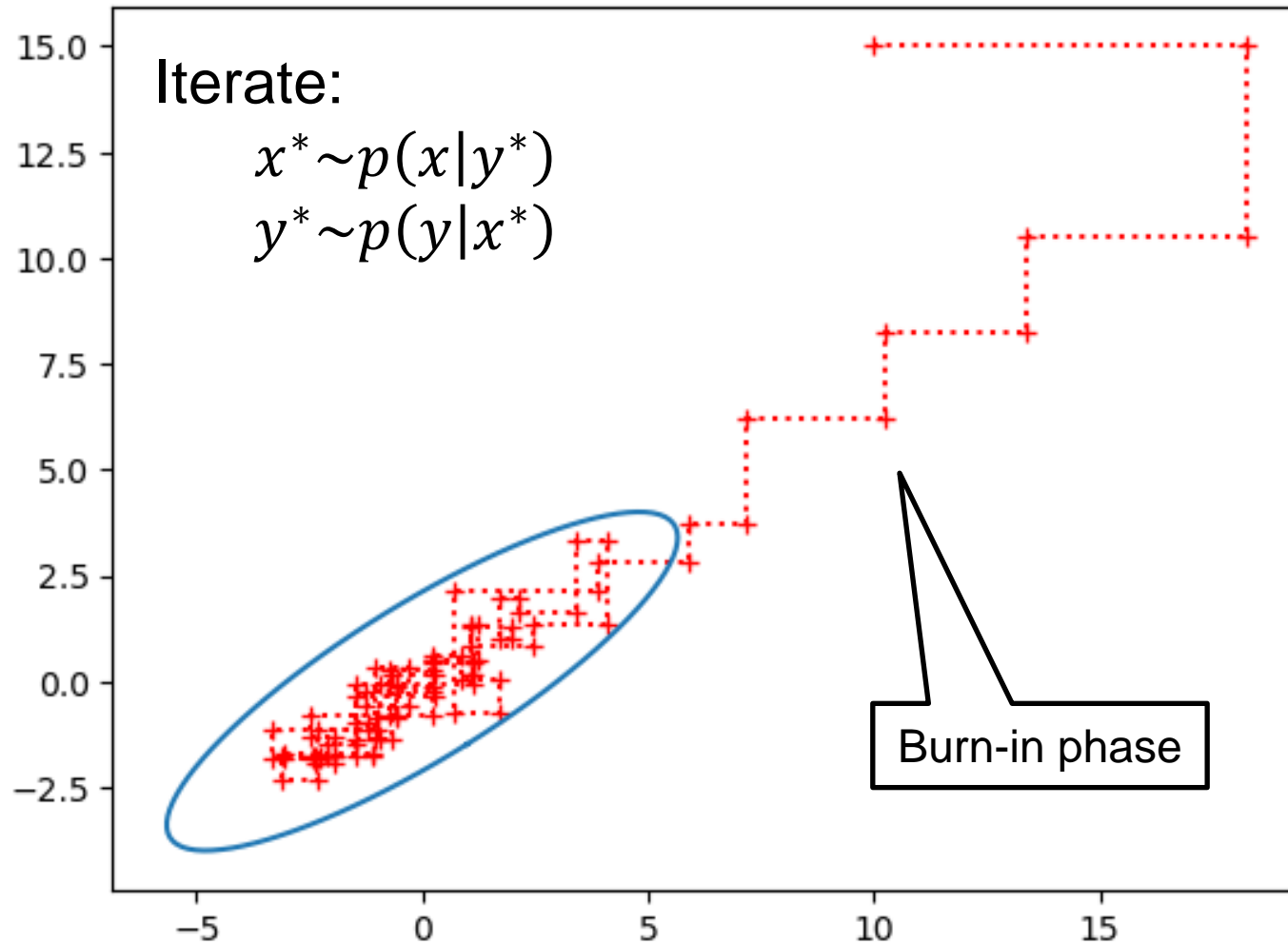


Gibbs Sampling

- Assume we cannot sample from the complex joint distribution $p(z_1, z_2)$ but it is possible to sample from the conditional distributions $p(z_1|z_2)$ and $p(z_2|z_1)$
 1. Initialize z_1^* to any value (i.e. chosen constant)
 2. Given current sample z_1^* generate $z_2^* \sim p(z_2|z_1)$
 3. Given current sample z_2^* generate $z_1^* \sim p(z_1|z_2)$
 4. Go to steps 2.
- In theory, after infinite number of iteration the final values z_1^*, z_2^* is a sample from $p(z_1, z_2)$
- Or, with increasing number of iterations, z_1^*, z_2^* converges to a valid sample from $p(z_1, z_2)$
- In practice, after several initial iterations (burn-in phase) take z_1^*, z_2^* from every N^{th} iteration and consider them samples from $p(z_1, z_2)$
 - Often $N = 1$ is used
 - Starting from a likely value of z_1^* requires less burn-in iterations
- This can be extended to any number variables
 - always sample one given current values for others
- Works for any random variables (discrete, continuous; scalars, vectors)

Gibbs sampling for 2D Gaussian

Of course, it is possible to efficiently and exactly sample directly from a 2D Gaussian distribution. We use this toy example only to demonstrate how Gibbs sampling works.



For 2D gaussian distribution

$$p\left(\begin{bmatrix} x \\ y \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} x \\ y \end{bmatrix} \middle| \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}\right)$$

the conditional probability

$$p(x|y) = \mathcal{N}(x | \mu_{x|y}, \Lambda_{xx}^{-1})$$

where

$$\mu_{x|y} = \mu_x - \Lambda_{xx}^{-1} \Lambda_{xy} (y - \mu_y)$$

and

$$\begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}^{-1} = \begin{bmatrix} \Lambda_{xx} & \Lambda_{xy} \\ \Lambda_{yx} & \Lambda_{yy} \end{bmatrix}$$

Approximate inference (for Bayesian LDA)

- Gibbs sampling
 - Instead of obtaining $p(\mathbf{Z}, \Theta, \Phi | \mathbf{W})$, we only generate samples from this distribution
 - We alternately sample from
 - » $\mathbf{Z}^* \sim P(\mathbf{Z} | \Theta, \Phi, \mathbf{W})$
 - » $\Theta^*, \Phi^* \sim P(\Theta, \Phi | \mathbf{Z}, \mathbf{W})$
- Variational Bayes
 - Approximate intractable $p(\mathbf{Z}, \Theta, \Phi | \mathbf{W})$ with tractable $q(\mathbf{Z}, \Theta, \Phi)$
 - Iteratively tune parameters of $q(\mathbf{Z}, \Theta, \Phi)$ to minimize $D_{KL}(q(\mathbf{Z}, \Theta, \Phi) || p(\mathbf{Z}, \Theta, \Phi | \mathbf{W}))$
- ...

GS inference for LDA: Sampling \mathbf{Z}

$$p(\mathbf{W}, \mathbf{Z}, \Theta, \Phi) = \left(\prod_{k=1}^K p(\boldsymbol{\varphi}_k) \right) \left(\prod_{d=1}^D p(\boldsymbol{\theta}_d) \right) \left(\prod_{d=1}^D \prod_{n=1}^{N_d} P(z_{dn} | \boldsymbol{\theta}_d) \right) \left(\prod_{d=1}^D \prod_{n=1}^{N_d} P(w_{dn} | \Phi, z_{dn}) \right)$$

Step 1:

- sample $\mathbf{Z}^* \sim P(\mathbf{Z} | \Theta, \Phi, \mathbf{W})$
- or equivalently for each d and n sample $z_{dn}^* \sim P(z_{dn} | \Theta, \Phi, \mathbf{W})$

We want posterior as a function of \mathbf{Z} , so we select only the terms involving \mathbf{Z} . The other terms are constant.

$$P(\mathbf{Z} | \Theta, \Phi, \mathbf{W}) \propto P(\mathbf{W}, \mathbf{Z}, \Theta, \Phi) \propto \prod_{d=1}^D \prod_{n=1}^{N_d} P(z_{dn} | \boldsymbol{\theta}_d) P(w_{dn} | \Phi, z_{dn}) \propto \prod_{d=1}^D \prod_{n=1}^{N_d} P(z_{dn} | \Theta, \Phi, \mathbf{W})$$

$$P(z_{dn} | \Theta, \Phi, \mathbf{W}) \propto P(w_{dn} | \Phi, z_{dn}) P(z_{dn} | \boldsymbol{\theta}_d) = \varphi_{z_{dn} w_{dn}} \theta_{d z_{dn}}$$

$$P(z_{dn} | \Theta, \Phi, \mathbf{W}) = \frac{\varphi_{z_{dn} w_{dn}} \theta_{d z_{dn}}}{\sum_{k'=1}^K \varphi_{k' w_{dn}} \theta_{d k'}} = \pi_{d w_{dn} z_{dn}}$$

$$\pi_{d v k} = \frac{\varphi_{k v} \theta_{d k}}{\sum_{k'=1}^K \varphi_{k' v} \theta_{d k'}}$$

Factorizes into product of independent terms one for each $z_{dn} \Rightarrow$ each z_{dn} can be sampled from independent $P(z_{dn} | \Theta, \Phi, \mathbf{W})$

GS inference for LDA: Sampling Θ and Φ

$$p(\mathbf{W}, \mathbf{Z}, \Theta, \Phi) \propto \left(\prod_{k=1}^K \prod_{v=1}^V \varphi_{kv}^{\beta_v - 1} \right) \left(\prod_{d=1}^D \prod_{k=1}^K \theta_{dk}^{\alpha_k - 1} \right) \left(\prod_{d=1}^D \prod_{k=1}^K \theta_{dk}^{c_{d,(\cdot)}^k} \right) \left(\prod_{k=1}^K \prod_{v=1}^V \varphi_{kv}^{c_{(\cdot),v}^k} \right)$$

Step 2:

- sample $\Theta^*, \Phi^* \sim P(\Theta, \Phi | \mathbf{Z}, \mathbf{W})$
- or equivalently for each $d = 1..D$ sample $\theta_d^* \sim P(\theta_d | \mathbf{Z}, \mathbf{W})$
and for each $k = 1..K$ sample $\varphi_k^* \sim P(\varphi_k | \mathbf{Z}, \mathbf{W})$

$$P(\Theta, \Phi | \mathbf{Z}, \mathbf{W}) \propto P(\mathbf{W}, \mathbf{Z}, \theta, \Phi) \propto \left(\prod_{d=1}^D \prod_{k=1}^K \theta_{dk}^{\alpha_k + c_{d,(\cdot)}^k - 1} \right) \left(\prod_{k=1}^K \prod_{v=1}^V \varphi_{kv}^{\beta_v + c_{(\cdot),v}^k - 1} \right)$$

$$P(\Theta, \Phi | \mathbf{Z}, \mathbf{W}) = \prod_{d=1}^D \text{Dir}(\theta_d | \alpha + \mathbf{C}_{d,(\cdot)}) \prod_{k=1}^K \text{Dir}(\varphi_k | \beta + \mathbf{C}_{(\cdot),k}^k) = \prod_{d=1}^D P(\theta_d | \mathbf{Z}, \mathbf{W}) \prod_{k=1}^K P(\varphi_k | \mathbf{Z}, \mathbf{W})$$

Factorizes into product of independent terms for each θ_d and $\varphi_k \Rightarrow$ can be sampled independently.

The dependence on \mathbf{Z} is through the counts $C_{d,v}^k$

where vectors $\mathbf{C}_{d,(\cdot)} = [C_{d,(\cdot)}^1, C_{d,(\cdot)}^2, \dots, C_{d,(\cdot)}^K]^T$ and $\mathbf{C}_{(\cdot),k}^k = [C_{(\cdot),k}^k, C_{(\cdot),k}^k, \dots, C_{(\cdot),k}^k]^T$

GS inference for LDA: Using word counts

- Until now, each training document $d = 1..D$ was represented by a variable length sequence of N_d words w_{dn} , where $n = 1..N_d$
- More conveniently, we can represent all documents by $D \times V$ matrix \mathbf{M} , with elements M_{dv} counting how many times document d contains word v
 - No need to know the order of the words w_{dn} in the sequence. The counts M_{dv} are enough.
 - Instead of sampling each z_{dn} , directly sample $C_{d,v}^k$

$$\pi_{dvk} = \frac{\varphi_{kv}\theta_{dk}}{\sum_{k'=1}^K \varphi_{k'v}\theta_{dk'}}$$

$$\boldsymbol{\pi}_{dv} = [\pi_{dv1}, \pi_{dv2}, \dots, \pi_{dvK}]^T$$

for $n = 1..N_d$

$$z_{dn}^* \sim P(z_{dn} | \boldsymbol{\Theta}, \boldsymbol{\Phi}, \mathbf{W}) = \text{Cat}(\boldsymbol{\pi}_{dw_{dn}})$$

$$\mathbf{C}_{d,v} = \begin{bmatrix} C_{d,v}^1 \\ C_{d,v}^2 \\ \vdots \\ C_{d,v}^K \end{bmatrix}$$

$$C_{d,v}^k = \sum_{n=1}^{N_d} \delta(z_{dn} = k) \delta(w_{dn} = v)$$



$$\mathbf{C}_{d,v} \sim \text{Multinomial}(\boldsymbol{\pi}_{dv}, M_{dv})$$

GS inference for LDA: Summary

for number of GS iterations

where

for $d = 1..D$

for $v = 1..V$

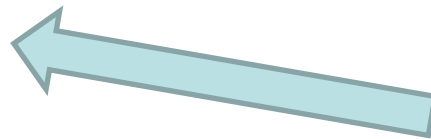
$$\mathbf{C}_{d,v} \sim \text{Multinomial}(\boldsymbol{\pi}_{dv}, M_{dv})$$



$$C_{d,(\cdot)}^k = \sum_{v=1}^V C_{d,v}^k, \quad C_{(\cdot),v}^k = \sum_{d=1}^D C_{d,v}^k$$

for $d = 1..D$

$$\boldsymbol{\theta}_d^* \sim P(\boldsymbol{\theta}_d | \mathbf{Z}, \mathbf{W}) = \text{Dir}(\boldsymbol{\theta}_d | \boldsymbol{\alpha} + \mathbf{C}_{d,(\cdot)})$$



$$\mathbf{C}_{d,(\cdot)} = \begin{bmatrix} C_{d,(\cdot)}^1 \\ C_{d,(\cdot)}^2 \\ \vdots \\ C_{d,(\cdot)}^K \end{bmatrix}, \quad \mathbf{C}_{(\cdot),v}^k = \begin{bmatrix} C_{(\cdot),v}^k \\ C_{(\cdot),v}^k \\ \vdots \\ C_{(\cdot),v}^k \end{bmatrix}$$

for $k = 1..K$

$$\boldsymbol{\varphi}_k^* \sim P(\boldsymbol{\varphi}_k | \mathbf{Z}, \mathbf{W}) = \text{Dir}(\boldsymbol{\varphi}_k | \boldsymbol{\beta} + \mathbf{C}_{(\cdot)}^k)$$

Where $\boldsymbol{\pi}_{dv}$ is evaluated using the currently values of $\boldsymbol{\Theta}$ and $\boldsymbol{\Phi}$ in each iteration.

After running a number of GS iterations, **we get likely samples** (from the posterior $p(\mathbf{Z}, \boldsymbol{\Theta}, \boldsymbol{\Phi} | \mathbf{W})$) of vectors $\boldsymbol{\theta}_d$ representing each document and $\boldsymbol{\varphi}_k$ representing latent topic distributions.

Approximate inference (for Bayesian LDA)

- Gibbs sampling
 - Instead of obtaining $p(\mathbf{Z}, \Theta, \Phi | \mathbf{W})$, we only generate samples from this distribution
- Variational Bayes
 - Approximate intractable $p(\mathbf{Z}, \Theta, \Phi | \mathbf{W})$ with tractable $q(\mathbf{Z}, \Theta, \Phi)$
 - Iteratively tune parameters of $q(\mathbf{Z}, \Theta, \Phi)$ to minimize $D_{KL}(q(\mathbf{Z}, \Theta, \Phi) || p(\mathbf{Z}, \Theta, \Phi | \mathbf{W}))$
- ...

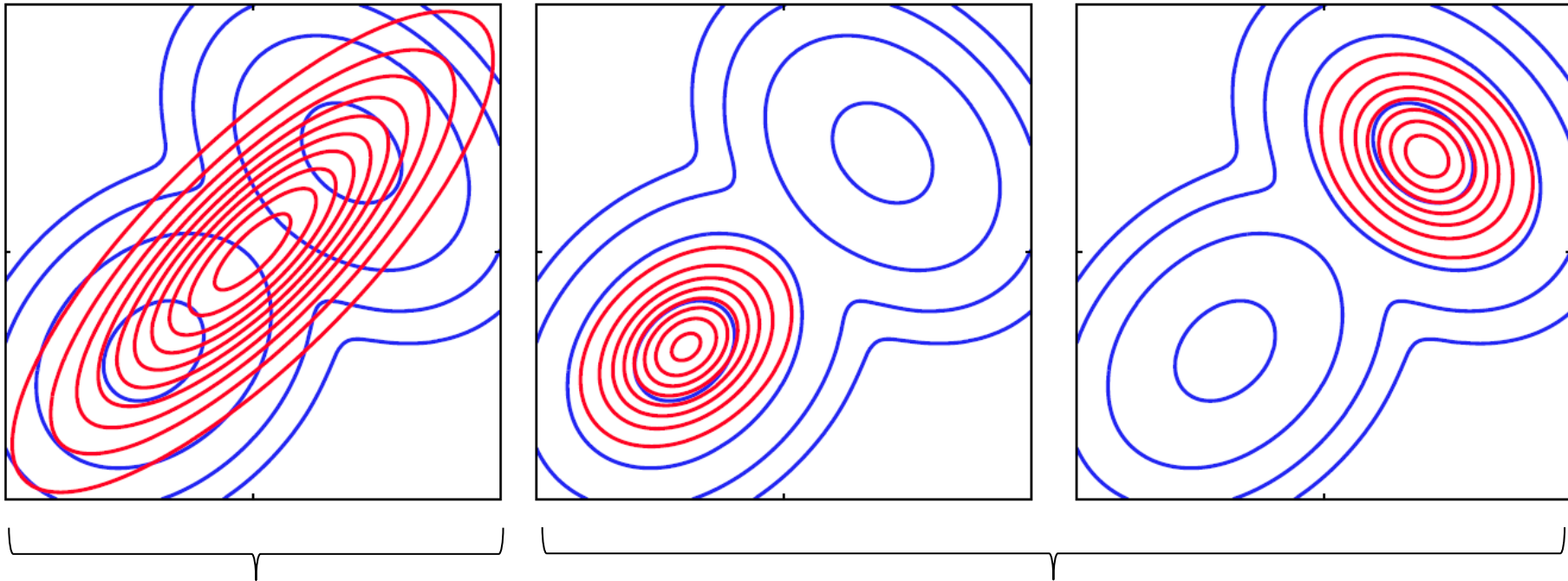
Variational Bayes

$$\ln p(\mathbf{X}) = \underbrace{\int q(\mathbf{Y}) \ln p(\mathbf{X}, \mathbf{Y}) \, d\mathbf{Y}}_{\mathcal{L}(q(\mathbf{Y}))} - \underbrace{\int q(\mathbf{Y}) \ln q(\mathbf{Y}) \, d\mathbf{Y}}_{D_{KL}(q(\mathbf{Y})||p(\mathbf{Y}|\mathbf{X}))} - \int q(\mathbf{Y}) \ln \frac{p(\mathbf{Y}|\mathbf{X})}{q(\mathbf{Y})} \, d\mathbf{Y}$$

- Find $q(\mathbf{Y})$, which is a good approximation for the true posterior $p(\mathbf{Y}|\mathbf{X})$
- Maximize $\mathcal{L}(q(\mathbf{Y}))$ w.r.t. $q(\mathbf{Y})$, which in turn minimizes $D_{KL}(q(\mathbf{Y})||p(\mathbf{Y}|\mathbf{X}))$
 - “Handcraft” a reasonable parametric distribution $q(\mathbf{Y}|\boldsymbol{\eta})$ and optimize $\mathcal{L}(q(\mathbf{Y}|\boldsymbol{\eta}))$ w.r.t. its parameters $\boldsymbol{\eta}$.
 - Mean field approximation assuming factorized form $q(\mathbf{Y})=q(\mathbf{Y}_1)q(\mathbf{Y}_2)q(\mathbf{Y}_3)\dots$

Minimizing Kullback-Leibler divergence

- We optimize parameters of (simpler) distribution $q(\mathbf{Y})$ to minimize Kullback-Leibler divergence between $q(\mathbf{Y})$ and $p(\mathbf{Y}|\mathbf{X})$.



- Minimizing $D_{KL}(p(\mathbf{Y}|\mathbf{X})||q(\mathbf{Y}))$.
- Not VB objective
- Expectation propagation
- Two local optima when (numerically) minimizing $D_{KL}(q(\mathbf{Y})||p(\mathbf{Y}|\mathbf{X}))$.
- VB performs this optimization

VB – Mean field approximation

- Popular Variational Bayes optimization method
- Variant of Variational Bayes, where the set of model variables \mathbf{Y} , can be split into subsets $\mathbf{Y}_1, \mathbf{Y}_2, \mathbf{Y}_3, \dots$, with **conditionally conjugate priors**
 - $p(\mathbf{Y}_i | \mathbf{X}, \mathbf{Y}_{\forall j \neq i})$ is tractable with conjugate prior
 - E.g. for Bayesian GMM $p(\mu_c, \lambda_c | \mathbf{X}, \mathbf{z})$ has NormalGamma prior
- We assume factorized approximate posterior

$$q(\mathbf{Y}) = q(\mathbf{Y}_1)q(\mathbf{Y}_2)q(\mathbf{Y}_3) \dots = \prod_i q(\mathbf{Y}_i)$$

- This factorization dictates the optimal (conjugate) distributions for the factors $q(\mathbf{Y}_i)$ and brings well defined iterative update formulas:

$$q(\mathbf{Y}_i)^* \propto \exp \left(\int q(\mathbf{Y}_{\forall j \neq i}) \ln p(\mathbf{X}, \mathbf{Y}) d\mathbf{Y}_{\forall j \neq i} \right)$$

Mean field - update

$$\begin{aligned}\mathcal{L}(q(\mathbf{Y})) &= \int q(\mathbf{Y}) \ln p(\mathbf{X}, \mathbf{Y}) \, d\mathbf{Y} - \int q(\mathbf{Y}) \ln q(\mathbf{Y}) \, d\mathbf{Y} = \int \prod_{i=1}^M q(\mathbf{Y}_i) \left[\ln p(\mathbf{X}, \mathbf{Y}) - \ln \prod_i q(\mathbf{Y}_i) \right] d\mathbf{Y} \\ &= \int \prod_{i=1}^M q(\mathbf{Y}_i) \left[\ln p(\mathbf{X}, \mathbf{Y}) - \sum_i \ln q(\mathbf{Y}_i) \right] d\mathbf{Y}\end{aligned}$$

- For example, let $M = 3$
- Now, let's optimize the lower bound $\mathcal{L}(q(\mathbf{Y}_1))$ w.r.t only one distribution $q(\mathbf{Y}_1)$

$$\begin{aligned}\mathcal{L}(q(\mathbf{Y}_1)) &= \iiint q(\mathbf{Y}_1)q(\mathbf{Y}_2)q(\mathbf{Y}_3) [\ln p(\mathbf{X}, \mathbf{Y}_1, \mathbf{Y}_2, \mathbf{Y}_3) - \ln q(\mathbf{Y}_1) - \ln q(\mathbf{Y}_2) - \ln q(\mathbf{Y}_3)] \, d\mathbf{Y}_1 \, d\mathbf{Y}_2 \, d\mathbf{Y}_3 \\ &= \int q(\mathbf{Y}_1) \underbrace{\iint q(\mathbf{Y}_2)q(\mathbf{Y}_3) \ln p(\mathbf{X}, \mathbf{Y}_1, \mathbf{Y}_2, \mathbf{Y}_3) \, d\mathbf{Y}_2 \, d\mathbf{Y}_3}_{\ln \tilde{p}(\mathbf{Y}_1) + const} \, d\mathbf{Y}_1 - \int q(\mathbf{Y}_1) \ln q(\mathbf{Y}_1) \, d\mathbf{Y}_1 + const \\ &= \int q(\mathbf{Y}_1) \ln \tilde{p}(\mathbf{Y}_1) \, d\mathbf{Y}_1 - \int q(\mathbf{Y}_1) \ln q(\mathbf{Y}_1) \, d\mathbf{Y}_1 + const = -D_{KL}(q(\mathbf{Y}_1) || \tilde{p}(\mathbf{Y}_1)) + const\end{aligned}$$

where $\tilde{p}(\mathbf{Y}_1)$ is normalized to be a valid distribution (therefore $+const$)

- $\mathcal{L}(q(\mathbf{Y}_1))$ is maximized by setting the D_{KL} term to zero, which implies

$$\ln q(\mathbf{Y}_1) = \ln \tilde{p}(\mathbf{Y}_1) = \iint q(\mathbf{Y}_2)q(\mathbf{Y}_3) \ln p(\mathbf{X}, \mathbf{Y}_1, \mathbf{Y}_2, \mathbf{Y}_3) \, d\mathbf{Y}_2 \, d\mathbf{Y}_3 + const$$

- In general, we can iteratively update each $q(\mathbf{Y}_i)$ given the others $q(\mathbf{Y}_{i \neq j})$ as:

$$q(\mathbf{Y}_j) \propto \exp \int q(\mathbf{Y}_{\forall j \neq i}) \ln p(\mathbf{X}, \mathbf{Y}) \, d\mathbf{Y}_{\forall j \neq i}$$

where each update guarantees to improve the lower bound $\mathcal{L}(q(\mathbf{Y}))$

Variational Bayes for LDA

Variational Bayes updates dictate:

$$q(\mathbf{Y}_i)^* \propto \exp \int q(\mathbf{Y}_{\forall j \neq i}) \ln p(\mathbf{X}, \mathbf{Y}) d\mathbf{Y}_{\forall j \neq i}$$

For the LDA model we chose to approximate the posterior using factorization

$$p(\mathbf{Z}, \Theta, \Phi | \mathbf{W}) \approx q(\mathbf{Z}, \Theta, \Phi) = q(\mathbf{Z})q(\Theta, \Phi)$$

Therefore, we search for updates in form:

$$q(\mathbf{Z})^* \propto \exp \iint q(\Theta, \Phi) \ln p(\mathbf{W}, \mathbf{Z}, \Theta, \Phi) d\Theta d\Phi$$

$$q(\Theta, \Phi)^* \propto \exp \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{W}, \mathbf{Z}, \Theta, \Phi)$$

Or equivalently:

$$\begin{aligned} \ln q(\mathbf{Z})^* &= \mathbb{E}_{q(\Theta, \Phi)} [\ln p(\mathbf{W}, \mathbf{Z}, \Theta, \Phi)] + \text{const} \\ \ln q(\Theta, \Phi)^* &= \mathbb{E}_{q(\mathbf{Z})} [\ln p(\mathbf{W}, \mathbf{Z}, \Theta, \Phi)] + \text{const} \end{aligned}$$

VB update for $q(\mathbf{Z})$

$$p(\mathbf{W}, \mathbf{Z}, \Theta, \Phi) = \left(\prod_{k=1}^K p(\boldsymbol{\varphi}_k) \right) \left(\prod_{d=1}^D p(\boldsymbol{\theta}_d) \right) \left(\prod_{d=1}^D \prod_{n=1}^{N_d} P(z_{dn} | \boldsymbol{\theta}_d) \right) \left(\prod_{d=1}^D \prod_{n=1}^{N_d} P(w_{dn} | \Phi, z_{dn}) \right)$$
$$\ln p(\mathbf{W}, \mathbf{Z}, \Theta, \Phi) = \sum_{k=1}^K \ln p(\boldsymbol{\varphi}_k) + \sum_{d=1}^D \left(\ln p(\boldsymbol{\theta}_d) + \sum_{n=1}^{N_d} \ln P(z_{dn} | \boldsymbol{\theta}_d) + \ln P(w_{dn} | \Phi, z_{dn}) \right)$$

The terms in $\ln p(\mathbf{W}, \mathbf{Z}, \Theta, \Phi)$ independent of z_{dn} are absorbed in const_2

$$\begin{aligned} \ln q(\mathbf{Z})^* &= \mathbb{E}_{q(\Theta, \Phi)} [\ln p(\mathbf{W}, \mathbf{Z}, \Theta, \Phi)] + \text{const}_1 \\ &= \sum_{d=1}^D \sum_{n=1}^{N_d} \mathbb{E}_{q(\Theta, \Phi)} [\ln P(z_{dn} | \boldsymbol{\theta}_d) + \ln P(w_{dn} | \Phi, z_{dn})] + \text{const}_2 \\ &= \sum_{d=1}^D \sum_{n=1}^{N_d} \ln q(z_{dn})^* \Rightarrow q(\mathbf{Z})^* = \prod_{d=1}^D \prod_{n=1}^{N_d} q(z_{dn})^* \end{aligned}$$

- We only require factorization $q(\mathbf{Z})q(\Theta, \Phi)$, but $q(\mathbf{Z})$ automatically further factorizes into a product of independent categorical distributions one for each z_{dn} - so called **induced factorization**
- We will derive the update for distributions $q(z_{dn})$ later.

VB update for $q(\Theta, \Phi)$

$$\ln p(\mathbf{W}, \mathbf{Z}, \Theta, \Phi) = \sum_{k=1}^K \ln p(\boldsymbol{\varphi}_k) + \sum_{d=1}^D \left(\ln p(\boldsymbol{\theta}_d) + \sum_{n=1}^{N_d} \ln \theta_{dz_{dn}} + \ln \varphi_{z_{dn} w_{dn}} \right) + \text{const}$$

$$\ln q(\Theta, \Phi)^* = \mathbb{E}_{q(\mathbf{Z})} [\ln p(\mathbf{W}, \mathbf{Z}, \Theta, \Phi)] + \text{const}$$

$$= \mathbb{E}_{q(\mathbf{Z})} \left[\sum_{d=1}^D \left(\ln p(\boldsymbol{\theta}_d) + \sum_{n=1}^{N_d} \ln \theta_{dz_{dn}} \right) + \left(\sum_{k=1}^K \ln p(\boldsymbol{\varphi}_k) + \sum_{d=1}^D \sum_{n=1}^{N_d} \ln \varphi_{z_{dn} w_{dn}} \right) \right] + \text{const}$$

$$= \sum_{d=1}^D \left(\ln p(\boldsymbol{\theta}_d) + \sum_{n=1}^{N_d} \mathbb{E}_{q(z_{dn})} [\ln \theta_{dz_{dn}}] \right) + \sum_{k=1}^K \left(\ln p(\boldsymbol{\varphi}_k) + \sum_{d=1}^D \sum_{n=1}^{N_d} \mathbb{E}_{q(z_{dn})} [\delta(z_{dn} = k) \ln \varphi_{k w_{dn}}] \right) + \text{const}$$

$$= \sum_{d=1}^D \ln q(\boldsymbol{\theta}_d)^* + \sum_{k=1}^K \ln q(\boldsymbol{\varphi}_k)^* \Rightarrow q(\Theta, \Phi)^* = \prod_{d=1}^D q(\boldsymbol{\theta}_d)^* \prod_{k=1}^K q(\boldsymbol{\varphi}_k)^*$$

induced factorization

VB update for $q(\boldsymbol{\theta}_d)$

$$\begin{aligned}
 \ln q(\boldsymbol{\theta}_d)^* &= \ln \text{Dir}(\boldsymbol{\theta}_d | \boldsymbol{\alpha}) + \sum_{n=1}^{N_d} \mathbb{E}_{q(z_{dn})} [\ln \theta_{dz_{dn}}] + \text{const} \\
 &= \left(\sum_{k=1}^K \alpha_k - 1 \ln \theta_{dk} \right) + \sum_{n=1}^{N_d} \sum_{k=1}^K q(z_{dn} = k) \ln \theta_{dk} + \text{const} \\
 &= \sum_{k=1}^K \left(\alpha_k - 1 + \sum_{n=1}^{N_d} q(z_{dn} = k) \right) \ln \theta_{dk} + \text{const} \\
 &= \sum_{k=1}^K \left(\alpha_k + \bar{C}_{d,(\cdot)}^k - 1 \right) \ln \theta_{dk} + \text{const} \\
 &= \ln \text{Dir}(\boldsymbol{\theta}_d | \boldsymbol{\alpha} + \bar{\mathbf{C}}_{d,(\cdot)})
 \end{aligned}$$

where

$$\bar{C}_{d,(\cdot)}^k = \sum_{n=1}^{N_d} q(z_{dn} = k)$$

$$\bar{\mathbf{C}}_{d,(\cdot)} = [\bar{C}_{d,(\cdot)}^1, \bar{C}_{d,(\cdot)}^2, \dots, \bar{C}_{d,(\cdot)}^K]$$

$$q(\boldsymbol{\theta}_d)^* = \text{Dir}(\boldsymbol{\theta}_d | \boldsymbol{\alpha} + \bar{\mathbf{C}}_{d,(\cdot)}) = \text{Dir}(\boldsymbol{\theta}_d | \boldsymbol{\alpha}_d^*)$$

Expected counts similar to the hard counts $C_{d,(\cdot)}^k$ from GS

Practically, the update means to calculate vector of parameters

$$\boldsymbol{\alpha}_d^* = \boldsymbol{\alpha} + \bar{\mathbf{C}}_{d,(\cdot)} \text{ for each } d$$

VB update for $q(\boldsymbol{\varphi}_k)$

$$\begin{aligned}
 \ln q(\boldsymbol{\varphi}_k)^* &= \ln \text{Dir}(\boldsymbol{\varphi}_k | \boldsymbol{\beta}) + \sum_{d=1}^D \sum_{n=1}^{N_d} \mathbb{E}_{q(z_{dn})} [\delta(z_{dn} = k) \ln \varphi_{kw_{dn}}] + \text{const} \\
 &= \left(\sum_{v=1}^V \beta_v - 1 \ln \varphi_{kv} \right) + \sum_{d=1}^D \sum_{n=1}^{N_d} q(z_{dn} = k) \ln \varphi_{kw_{dn}} + \text{const} \\
 &= \sum_{v=1}^V \left(\beta_v - 1 + \sum_{d=1}^D \sum_{n=1}^{N_d} q(z_{dn} = k) \delta(w_{dn} = v) \right) \ln \varphi_{kv} + \text{const} \\
 &= \sum_{v=1}^V (\beta_v + \bar{c}_{(\cdot),v}^k - 1) \ln \varphi_{kv} + \text{const} \\
 &= \ln \text{Dir}(\boldsymbol{\varphi}_k | \boldsymbol{\beta} + \bar{\mathbf{c}}_{(\cdot)}^k)
 \end{aligned}$$

where

$$\begin{aligned}
 \bar{c}_{(\cdot),v}^k &= \sum_{d=1}^D \sum_{n=1}^{N_d} q(z_{dn} = k) \delta(w_{dn} = v) \\
 \bar{\mathbf{c}}_{(\cdot)}^k &= [\bar{c}_{(\cdot),1}^k, \bar{c}_{(\cdot),2}^k, \dots, \bar{c}_{(\cdot),V}^k]
 \end{aligned}$$

$$q(\boldsymbol{\varphi}_k)^* = \text{Dir}(\boldsymbol{\varphi}_k | \boldsymbol{\beta} + \bar{\mathbf{c}}_{(\cdot)}^k) = \text{Dir}(\boldsymbol{\varphi}_k | \boldsymbol{\beta}_k^*)$$

Practically, the update means to calculate vector of parameters

$$\boldsymbol{\beta}_k^* = \boldsymbol{\beta} + \bar{\mathbf{c}}_{(\cdot)}^k$$

for each k

VB update for $q(z_{dn})$

$$\begin{aligned}\ln q(z_{dn})^* &= \mathbb{E}_{q(\boldsymbol{\theta}, \boldsymbol{\Phi})} [\ln P(z_{dn} | \boldsymbol{\theta}_d) + \ln P(w_{dn} | \boldsymbol{\Phi}, z_{dn})] + \text{const} \\ &= \mathbb{E}_{q(\boldsymbol{\theta}_d)} [\ln P(z_{dn} | \boldsymbol{\theta}_d)] + \mathbb{E}_{q(\boldsymbol{\varphi}_k)} [\ln P(w_{dn} | \boldsymbol{\Phi}, z_{dn})] + \text{const} \\ &= \mathbb{E}_{q(\boldsymbol{\theta}_d)} [\ln \theta_{dz_{dn}}] + \mathbb{E}_{q(\boldsymbol{\varphi}_k)} [\ln \varphi_{z_{dn} w_{dn}}] + \text{const}\end{aligned}$$

For $q(\boldsymbol{\theta}_d)^* = \text{Dir}(\boldsymbol{\theta}_d | \boldsymbol{\alpha}_d^*)$ and $q(\boldsymbol{\varphi}_k)^* = \text{Dir}(\boldsymbol{\varphi}_k | \boldsymbol{\beta}_k^*)$

$$\begin{aligned}\ln \rho_{dvk} &= \mathbb{E}_{q(\boldsymbol{\theta}_d)} [\ln \theta_{dk}] + \mathbb{E}_{q(\boldsymbol{\varphi}_k)} [\ln \varphi_{kv}] + \text{const} \\ &= \psi(\alpha_{dk}^*) - \psi\left(\sum_{k'=1}^K \alpha_{dk'}^*\right) + \psi(\beta_{kv}^*) - \psi\left(\sum_{v'=1}^V \beta_{kv'}^*\right) + \text{const}\end{aligned}$$

Therefore, we update

$$q(z_{dn} = k)^* = \pi_{dw_{dn}k}$$

$$\pi_{dvk} = \frac{\rho_{dvk}}{\sum_{k'=1}^K \rho_{dvk'}}$$

Practically, the update means to evaluate 3D matrix with elements π_{dvk}

VB inference for LDA using word counts

- Again, as for the GS inference, rather than using sequence of words w_{dn} , we prefer to represent documents by $D \times V$ matrix M , with elements M_{dv} counting how many times document d contains word v
- Instead of estimating each $q(z_{dn} = k)$, we directly estimate expected counts $\bar{C}_{(\cdot),v}^k$ and $\bar{C}_{d,(\cdot)}^k$

$$q(z_{dn} = k) = \pi_{dw_{dn}k}$$

$$\bar{C}_{(\cdot),v}^k = \sum_{d=1}^D \bar{C}_{d,v}^k = \sum_{d=1}^D \sum_{n=1}^{N_d} q(z_{dn} = k) \delta(w_{dn} = v)$$

$$\bar{C}_{d,(\cdot)}^k = \sum_{v=1}^V \bar{C}_{d,v}^k = \sum_{n=1}^{N_d} q(z_{dn} = k)$$

$$\bar{C}_{d,v}^k = \sum_{n=1}^{N_d} q(z_{dn} = k) \delta(w_{dn} = v) = \pi_{dvk} \sum_{n=1}^{N_d} \delta(w_{dn} = v) \quad \longrightarrow \quad \bar{C}_{d,v}^k = \pi_{dvk} M_{dv}$$

ELBO objective to monitor progress

$$\begin{aligned}
 \mathcal{L}(q(\mathbf{Z}, \boldsymbol{\Theta}, \boldsymbol{\Phi})) &= \mathbb{E}_{q(\mathbf{Z}, \boldsymbol{\Theta}, \boldsymbol{\Phi})} \left[\ln \frac{p(\mathbf{W}, \mathbf{Z}, \boldsymbol{\Theta}, \boldsymbol{\Phi})}{q(\mathbf{Z}, \boldsymbol{\Theta}, \boldsymbol{\Phi})} \right] \\
 &= \sum_{d=1}^D \sum_{n=1}^{N_d} \mathbb{E}_{q(z_{dn})} \left[\mathbb{E}_{q(\boldsymbol{\theta}_d)} [\ln P(z_{dn} | \boldsymbol{\theta}_d)] + \mathbb{E}_{q(\boldsymbol{\varphi}_k)} [\ln P(w_{dn} | \boldsymbol{\Phi}, z_{dn})] - \ln q(z_{dn}) \right] \\
 &\quad - \sum_{k=1}^K \mathbb{E}_{q(\boldsymbol{\varphi}_k)} \left[\ln \frac{q(\boldsymbol{\varphi}_k)}{p(\boldsymbol{\varphi}_k)} \right] - \sum_{d=1}^D \mathbb{E}_{q(\boldsymbol{\theta}_d)} \left[\ln \frac{q(\boldsymbol{\theta}_d)}{p(\boldsymbol{\theta}_d)} \right] \\
 &= \sum_{d=1}^D \sum_{n=1}^{N_d} \sum_{k=1}^K q(z_{dn} = k) [\ln \rho_{dw_{dn}k} - \ln q(z_{dn} = k)] - \sum_{d=1}^D KL(q(\boldsymbol{\theta}_d) || p(\boldsymbol{\theta}_d)) - \sum_{k=1}^K KL(q(\boldsymbol{\varphi}_k) || p(\boldsymbol{\varphi}_k)) \\
 &= \sum_{d=1}^D \sum_{v=1}^V \sum_{k=1}^K \bar{C}_{d,v}^k (\ln \rho_{dvk} - \ln \pi_{dvk}) - \sum_{d=1}^D KL(\text{Dir}(\boldsymbol{\alpha}_d^*) || \text{Dir}(\boldsymbol{\alpha})) - \sum_{k=1}^K KL(\text{Dir}(\boldsymbol{\beta}_k^*) || \text{Dir}(\boldsymbol{\beta}))
 \end{aligned}$$

$$KL(\text{Dir}(\boldsymbol{\alpha}) || \text{Dir}(\boldsymbol{\beta})) = \ln \Gamma\left(\sum_{c=1}^C \alpha_c\right) - \sum_{c=1}^C \ln \Gamma(\alpha_c) - \ln \Gamma\left(\sum_{c=1}^C \beta_c\right) + \sum_{c=1}^C \ln \Gamma(\beta_c) + \sum_{c=1}^C (\alpha_c - \beta_c) \left(\psi(\alpha_c) - \psi\left(\sum_{c'=1}^C \alpha_{c'}\right) \right)$$

Efficient ELBO calculation

$$\mathcal{L}(q(\mathbf{Z}, \Theta, \Phi)) = \sum_{d=1}^D \sum_{v=1}^V \sum_{k=1}^K \bar{C}_{d,v}^k (\ln \rho_{dvk} - \ln \pi_{dvk}) - \sum_{d=1}^D KL(\text{Dir}(\alpha_d^*) || \text{Dir}(\alpha)) - \sum_{k=1}^K KL(\text{Dir}(\beta_k^*) || \text{Dir}(\beta))$$

Right after updating $q(z_{dn})$ (i.e. evaluating the terms π_{dvk}), the **red term** becomes independent of k since

$$\pi_{dvk} = \frac{\rho_{dvk}}{\sum_{k'=1}^K \rho_{dvk'}} \Rightarrow \ln \rho_{dvk} - \ln \pi_{dvk} = \ln \sum_{k'=1}^K \rho_{dvk'}$$

Therefore, right after updating $q(z_{dn})$ and before updating any $q(\theta_d)$ or $q(\varphi_k)$, we efficiently calculate ELBO $\mathcal{L}(q(\mathbf{Z}, \Theta, \Phi))$ as:

$$\mathcal{L}(q(\mathbf{Z}, \Theta, \Phi)) = \sum_{d=1}^D \sum_{v=1}^V M_{dv} \ln \sum_{k=1}^K \rho_{dvk} - \sum_{d=1}^D KL(\text{Dir}(\alpha_d^*) || \text{Dir}(\alpha)) - \sum_{k=1}^K KL(\text{Dir}(\beta_k^*) || \text{Dir}(\beta))$$

where we have used $\sum_{k=1}^K \bar{C}_{d,v}^k = M_{dv}$, and where we can reuse the terms $\sum_{k'=1}^K \rho_{dvk'}$ that were just calculated for normalizing the terms π_{dvk} .

VB inference for LDA: Summary

for number of VB iterations

for every $d = 1..D$, $v = 1..V$ and $k = 1..K$

$$\ln \rho_{dvk} = \psi(\alpha_{dk}^*) - \psi\left(\sum_{k'=1}^K \alpha_{dk'}^*\right) + \psi(\beta_{kv}^*) - \psi\left(\sum_{v'=1}^V \beta_{kv'}^*\right)$$

$$\pi_{dvk} = \frac{\rho_{dvk}}{\sum_{k'=1}^K \rho_{dvk'}}$$

$$\bar{C}_{d,v}^k = \pi_{dvk} M_{dv}$$

for $d = 1..D$

$$\alpha_d^* = \alpha + \bar{\mathbf{C}}_{d,(\cdot)}$$

for $k = 1..K$

$$\beta_k^* = \beta + \bar{\mathbf{C}}_{(\cdot)}^k$$

where

$$\bar{C}_{d,(\cdot)}^k = \sum_{v=1}^V \bar{C}_{d,v}^k, \quad \bar{C}_{(\cdot),v}^k = \sum_{d=1}^D \bar{C}_{d,v}^k$$

$$\bar{\mathbf{C}}_{d,(\cdot)} = \begin{bmatrix} C_{d,(\cdot)}^1 \\ C_{d,(\cdot)}^2 \\ \vdots \\ C_{d,(\cdot)}^K \end{bmatrix}, \quad \bar{\mathbf{C}}_{(\cdot)}^k = \begin{bmatrix} C_{(\cdot),1}^k \\ C_{(\cdot),2}^k \\ \vdots \\ C_{(\cdot),V}^k \end{bmatrix}$$

Where π_{dv} is evaluated using the currently values of α_d^* and β_k^* in each iteration.

After running a number of VB iterations, $q(\boldsymbol{\theta}_d)^* = \text{Dir}(\boldsymbol{\theta}_d | \alpha_d^*)$ and $q(\boldsymbol{\varphi}_k)^* = \text{Dir}(\boldsymbol{\varphi}_k | \beta_k^*)$ are approximate posteriors for all $d = 1..D$ and $k = 1..K$